

On Shore and Johnson properties for a Special Case of Csiszár f -divergences

Jithin Vachery and Ambedkar Dukkipati
 Dept of CSA, IISc
 {jithinvachery, ambedkar} @csa.iisc.ernet.in

Abstract—The importance of power-law distributions is attributed to the fact that most of the naturally occurring phenomenon exhibit this distribution. While exponential distributions can be derived by minimizing KL-divergence w.r.t some moment constraints, some power law distributions can be derived by minimizing some generalizations of KL-divergence (more specifically some special cases of Csiszár f -divergences). Divergence minimization is very well studied in information theoretical approaches to statistics. In this work we study properties of minimization of Tsallis divergence, which is a special case of Csiszár f -divergence. In line with the work by Shore and Johnson (IEEE Trans. IT, 1981), we examine the properties exhibited by these minimization methods including the Pythagorean property.

I. INTRODUCTION

Shannon measure of information, also called entropy, is central to information theory which has wide range of applications spanning, communication theory, statistical mechanics, probability theory, statistical inference etc. [1]. It quantifies uncertainty or information that is associated with a discrete random variable by taking an average of uncertainty (Hartley information) associated with each state. The first generalization of this measure of information was suggested by Rényi [2]. He replaced the linear averaging by K-N averages (Kolmogorov-Nagumo averages) and imposed additivity constraint. Havrda and Charvat [3] introduced one more generalization which is now known as nonextensive entropy or Tsallis entropy [4], [5], [6], which has been studied in statistical mechanics. Another important notion is that of finding the distance or divergence between two probability distributions. The information measure capturing this is KL-divergence, which is the directed distance between two probability distributions. KL-divergence is a special case of Tsallis divergence, which in turn is a special case of Csiszár f -divergence [7]. KL-divergence plays a central role in Kullback's minimum divergence principle, Which is a means of estimating the probability distribution of a system. It suggests the minimization of KL-divergence using a given prior distribution, subject to moment constraints as the estimation technique. Kullback's minimum divergence principle reduces to Jaynes maximum entropy principle

when we use uniform distribution as the prior. Kullback's minimum divergence principle can be extended to generalized divergences. When applied to classical KL-divergence, this yields a distribution from the exponential family. Whereas applying Kullback's principle to Tsallis divergence gives a power-law distribution.

Exponential distributions are very important class of distributions and many problems have been successfully modeled using this [8]. Though exponential distributions are used in many modeling problems [9] due to theoretical tractability, many naturally occurring phenomena exhibit power-law distributions. It is of great practical and theoretical interest to study both these family of distributions. In this work we have been able to establish many properties for Tsallis divergence. We have established the property of *transformation invariance* and *subset independence*. In addition we have found some properties for Tsallis divergence minimization in classical constraints viz. *uniqueness*, *reflexiveness*, *idempotence*, *invariance*, *weak subset independence* and *subset aggregation*. In this work we have also attempted to derive a *Pythagorean property*. In addition we have proposed a $q \leftrightarrow 2 - q$ *additive transformation* for Tsallis divergence.

The paper is organized as follows. In Section II we introduce the preliminaries and basics required for understanding the results. Sections III through V are dedicated to the results and observations made. In these sections we perform Tsallis divergence minimization for classical constraints and we follow it up with the analysis of the properties exhibited by the same. In particular we are study about the Shore and Johnson properties. In the subsequent section we discuss about the a transformation relation which we established.

II. PRELIMINARIES AND BACKGROUND

A. Exponential family and KL Divergence

In many of the problems we might have a prior estimate of the probability distribution and given such a prior we are interested in finding the probability distribution that is closest to this prior, which also satisfies the set of linear constraints. To define the notion of closeness we

need a distance measure between two distributions. One such distance measure is KL divergence [10] defined as

$$I(p||r) = \sum_{x \in \mathcal{X}} p(x) \ln \left(\frac{p(x)}{r(x)} \right) ,$$

where r is the prior. The minimization of KL-divergence results in a posterior which is from the exponential family.

B. Power-Law distribution and Generalized Divergence

f -divergence is a generalized measure of divergence, that was introduced by Csiszár [7] and independently by Ali & Silvey [11]. Let $f(t)$ be a real valued convex function defined for $t > 0$, with $f(1) = 0$. The f -divergence of a distribution p from r is defined by

$$D_f(p||r) = \sum_{x \in \mathcal{X}} r(x) f \left(\frac{p(x)}{r(x)} \right) .$$

Here we take $0f(\frac{0}{0}) = 0, f(0) = \lim_{t \rightarrow 0} f(t)$. f -divergence has many important properties like non-negativity, monotonicity and convexity. This has been used in many applications like speech recognition [12], analysis of contingency tables [7], etc. By specializing f to various functions we get different divergences like KL-divergence, χ^2 -divergence, Hellinger distance, variational distance, Tsallis-divergence, etc. On setting $f(t) = t \ln_q t$ we get Tsallis divergence [4], defined as

$$I_q(p||r) = - \sum_{x \in \mathcal{X}} p(x) \ln_q \frac{r(x)}{p(x)} ,$$

where \ln_q is q -logarithm function [13], defined as, $\ln_q x = \frac{x^{(1-q)} - 1}{1-q}$ ($x > 0, q \in \mathbb{R}$). Tsallis divergence recovers KL-divergence for $q \rightarrow 1$ i.e., $\lim_{q \rightarrow 1} I_q(p||r) = I(p||r)$. For values of $q > 0$ we have $I_q(p||r) \geq 0$ and Tsallis divergence becomes a convex function of both the parameters. Tsallis divergence also exhibits pseudo additivity property, i.e., $I_q(X1 \times X2 || Y1 \times Y2) = I_q(X1 || X2) \oplus_q I_q(Y1 || Y2)$, where $X1$ and $X2$ are independent, so are $Y1$ and $Y2$. Here \oplus_q is addition in q -deformed algebra [13] defined as, $x \oplus_q y = x + y + (1-q)xy$. In the minimization of Tsallis divergence the choice of constraints play an important role [14].

Tsallis Divergence minimization with respect to q -expectation constraint has been studied by [15]. In this case Pythagoras theorem is established by [16], [17], [18] and proved in differential geometric setup by Ohara [19]. Tsallis divergence minimization with normalized constraints gives probability distribution which is self referential in nature, i.e., $p(x)$ depends of $p(x)$. Here too we have nonextensive Pythagoras property [16], [17] exhibited by Tsallis-divergence.

In this paper we are going to study this minimization with respect to classical expectations, as it has the important property of convexity, ensuring a unique solution.

III. BASIC SHORE AND JOHNSON PROPERTIES

Shore and Johnson [20] in their work in 1981 had discussed many of the important properties of KL-divergence minimization. We have found that many of those properties hold in the case of Tsallis divergence. In this section we shall discuss about the properties that pertain to Tsallis divergence, i.e., regardless of minimization.

In this section and section V we shall be using the following notation.

Let p be a pmf. on random variable X taking values from \mathcal{X} . We would like to impose the following linear equality and inequality constraints on it.

$$\sum_{x \in \mathcal{X}} p(x) = 1 , \quad (1)$$

$$\sum_{x \in \mathcal{X}} p(x) u_m = \langle u_m \rangle \quad m = 1 \dots M , \quad (2)$$

$$\sum_{x \in \mathcal{X}} p(x) w_n \geq \langle w_n \rangle \quad n = 1 \dots N . \quad (3)$$

Equations (1),(2) and (3) constitute the constraint set. This can also be considered as the information available about the probability distribution. We shall denote a constraint set by \mathcal{C} , and a subscript to distinguish between different constraint sets.

Hence the task of divergence minimization can be viewed as, given a prior probability distribution $q(x)$ and constraint set \mathcal{C} finding the probability distribution p_{min} such that $p_{min} = \arg \min_{p \in \mathcal{C}} I_q(p||r)$. It can be easily verified that the constraint set \mathcal{C} constitutes a convex set. We would like to inform that some of these notation have been borrowed from [20].

Invariance of KL-divergence to coordinate transformations enables us to generalize KL-divergence to continuous random variables. We have observed that the invariance property holds true in the case of Tsallis divergence too.

Proposition 1 (Invariance): Let Γ be a coordinate transformation from $x \in \mathcal{X}$ to $y \in \mathcal{X}'$ with $(\Gamma p)(y) = J^{-1}p(x)$, where J is the Jacobian $J = \partial(y)/\partial(x)$. Let $\Gamma \mathcal{X}$ be the set of densities Γp corresponding to densities $p \in \mathcal{X}$. Let $(\Gamma \mathcal{C}) \subseteq (\Gamma \mathcal{X})$ correspond to $\mathcal{C} \subseteq \mathcal{X}$. Then, given a prior distribution r

$$\arg \min_{p \in \Gamma \mathcal{C}} I_q(p||\Gamma r) = \arg \min_{s \in \mathcal{C}} I_q(s||r) , \quad (4)$$

$$\text{and } I_q(\Gamma p_{min}||\Gamma r) = I_q(p_{min}||r) , \quad (5)$$

hold. where $\Gamma p_{min} = \arg \min_{p \in \Gamma \mathcal{C}} I_q(p||\Gamma r)$ and $p_{min} = \arg \min_{s \in \mathcal{C}} I_q(s||r)$.

Proof: We have $(\Gamma p)(y) = J^{-1}p(x)$, where J is the Jacobian $J = \partial(y)/\partial(x)$.

$$\begin{aligned} I_q(\Gamma p || \Gamma r) &= - \int_{\Gamma \mathcal{X}} \Gamma p(y) \ln_q \left(\frac{\Gamma r(y)}{\Gamma p(y)} \right) dy \\ &= - \int_{\mathcal{X}} J^{-1}p(x) \ln_q \left(\frac{J^{-1}r(x)}{J^{-1}p(x)} \right) J dx \\ &= - \int_{\mathcal{X}} p(x) \ln_q \frac{r(x)}{p(x)} dx \\ &= I_q(p || r) \end{aligned}$$

This proves (5). From (5) it also follows that the minimum in $\Gamma \mathcal{C}$ corresponds to the minimum in \mathcal{C} , which proves (4). \blacksquare

Proposition 2 (Subset Independence): Let S_1, S_2, \dots, S_n be a partition of \mathcal{X} . Let the new information \mathcal{C} comprise about each of the conditional densities $p(x/x \in S_i)$, $i = 1 \dots n$. Thus, $\mathcal{C} = \mathcal{C}_1 \wedge \mathcal{C}_2 \wedge \dots \wedge \mathcal{C}_n$, where \mathcal{C}_i is the constraint set on the conditional densities of S_i . Let \mathcal{M} be the new information giving the probability of being in each of the n subsets, which is the constraint

$$\sum_{x \in S_i} p(x) = m_i, \quad i = 1 \dots n,$$

where m_i are known values. Then given the prior distribution r ,

$$p_{\mathcal{C}\mathcal{M}}^{min}(x/x \in S_i) = \arg \min_{p \in \mathcal{C}_i} I_q(p || r_i), \quad q \in (0, 1), \quad (6)$$

and

$$\begin{aligned} I_q(p_{\mathcal{C}\mathcal{M}}^{min} || r) &= \sum_{i=1}^n m_i I_q(p_i || r_i) - \sum_{i=1}^n m_i \ln_q \frac{S_i}{m_i} \\ &+ (1-q) \sum_{i=1}^n \left(m_i \ln_q \frac{S_i}{m_i} I_q(p_i || r_i) \right) \end{aligned} \quad (7)$$

hold, where

$$\begin{aligned} p_{\mathcal{C}\mathcal{M}}^{min} &= \arg \min_{p \in \mathcal{C} \wedge \mathcal{M}} I_q(p || r), \\ p_i(x) &= p_{\mathcal{C}\mathcal{M}}^{min}(x/x \in S_i), \\ r_i(x) &= r(x/x \in S_i), \end{aligned}$$

and s_i are the prior probability of being in each subset, given by $s_i = \sum_{x \in S_i} r(x)$.

Proof:

$$I_q(p_{\mathcal{C}\mathcal{M}}^{min} || r) = - \sum_{i=1}^n \sum_{x \in S_i} m_i p_i(x) \ln_q \frac{s_i r_i(x)}{m_i r_i(x)}.$$

Using the relation $\ln_q(xy) = \ln_q x + \ln_q y + (1-q) \ln_q x \ln_q y$, we get

$$\begin{aligned} I_q(p_{\mathcal{C}\mathcal{M}}^{min} || r) &= - \sum_{i=1}^n \sum_{x \in S_i} m_i p_i(x) \left(\ln_q \frac{s_i}{m_i} + \ln_q \frac{r_i(x)}{p_i(x)} \right. \\ &\quad \left. + (1-q) \ln_q \frac{s_i}{m_i} \ln_q \frac{r_i(x)}{p_i(x)} \right) \\ &= \sum_{i=1}^n m_i I_q(p_i || r_i) - \sum_{i=1}^n m_i \ln_q \frac{s_i}{m_i} \\ &\quad + (1-q) \sum_{i=1}^n \left(m_i \ln_q \frac{s_i}{m_i} I_q(p_i || r_i) \right), \end{aligned}$$

this proves (7). To prove (6) it may be noted that each of the terms $m_i \ln_q \frac{s_i}{m_i}$ is a constant. Hence minimizing rhs of (7) is independent of the values taken by it. i.e for $q \in (0, 1)$ minimizing $I_q(p_{\mathcal{C}\mathcal{M}}^{min} || r)$ is equivalent to minimizing each of the terms, $I_q(p_i || r_i)$. \blacksquare

Let us further analyze equation (6) and try to interpret it. What this means is that, given a system which naturally partitions into subsets, we can find the posterior densities in two different ways

- 1) We can find the posterior $p_{\mathcal{C}\mathcal{M}}^{min}$ and condition it on the different subsets S_i or
- 2) We can condition the prior r on the different subsets S_i and use that as a prior to minimize in the constraint set \mathcal{C}_i

By (6) both these approaches should give the same result.

IV. TSALLIS DIVERGENCE MINIMIZATION - CLASSICAL

The task of minimization can be defined as follows: Minimize $I_q(p || r)$ subject to the constraints

$$\begin{aligned} \sum_{x \in \mathcal{X}} p(x) &= 1, \\ p(x) &\geq 0, \\ \sum_{x \in \mathcal{X}} u_m(x) p(x) &= \langle u_m \rangle, \quad m = 1, \dots, M. \end{aligned} \quad (8)$$

By choosing the Lagrangian for the minimization problem as

$$\begin{aligned} \mathcal{L} &= \sum_{x \in \mathcal{X}} p(x) \frac{\left[\frac{p(x)}{r(x)} \right]^{q-1} - 1}{q-1} - \left(\frac{q\lambda-1}{q-1} \right) (\sum_{x \in \mathcal{X}} p(x) - 1) \\ &\quad - \sum_{m=1}^M q\lambda\beta_m (\sum_{x \in \mathcal{X}} u_m(x) p(x) - \langle u_m \rangle). \end{aligned}$$

The distribution that we get after minimization is

$$p(x) = r(x) \left[\lambda \left(1 + (q-1) \sum_{m=1}^M \beta_m u_m(x) \right) \right]^{\frac{1}{q-1}}. \quad (9)$$

Substituting (9) in (8) we get

$$\lambda^{\frac{1}{q-1}} = \frac{1}{\sum_{x \in \mathcal{X}} \left[r(x) \left(1 - (1-q) \sum_{m=1}^M \beta_m u_m(x) \right)^{\frac{1}{q-1}} \right]} .$$

Substituting in (9) we get

$$p(x) = \frac{r(x) \left(1 - (1-q) \sum_{m=1}^M \beta_m u_m(x) \right)^{\frac{1}{q-1}}}{\widehat{Z}} , \quad (10)$$

where

$$\widehat{Z} = \sum_{x \in \mathcal{X}} \left[r(x) \left(1 - (1-q) \sum_{m=1}^M \beta_m u_m(x) \right)^{\frac{1}{q-1}} \right] .$$

equation (10) can be rewritten as

$$p(x) = \frac{r(x)}{\widehat{Z} \exp_q \left(- \sum_{m=1}^M \beta_m u_m(x) \right)} . \quad (11)$$

Where \exp_q is exponentiation in q -deformed algebra [13], and is defined as,

$$\exp_q(x) = \begin{cases} [1 + (1-q)x]^{\frac{1}{1-q}} & \text{if } 1 + (1-q)x \geq 0 \\ 0 & \text{otherwise} . \end{cases}$$

using the relation $\frac{1}{\exp_q(x)} = \exp_q \left(\frac{-x}{1+(1-q)x} \right)$, we get

$$p(x) = \frac{r(x)}{\widehat{Z}} \exp_q \left(\frac{\sum_{m=1}^M \beta_m u_m(x)}{1 - (1-q) \sum_{m=1}^M \beta_m u_m(x)} \right) . \quad (12)$$

Note that we need an extra condition known as *Tsallis cut-off condition* to prevent negative values for $p(x)$. We have assumed this condition to be implicit.

V. SHORE AND JOHNSON PROPERTIES INVOLVING MAXIMUM ENTROPY

In this section we shall discuss properties which depend on the formalism employed.

Proposition 3 (Uniqueness): For $q > 0$ given a prior, the posterior probability distribution is unique.

Proof: For $q > 0$ Tsallis divergence is a convex function, for both its parameter. Since the constraint set \mathcal{C} is a convex set, the minimization is always unique. ■

Proposition 4 (Reflexiveness): For $q > 0$, given a prior r and constraint set \mathcal{C} , the posterior obtained by minimizing the Tsallis divergence is same as r if and only if $r \in \mathcal{C}$

Proof: This property follows directly from the following facts $I_q(p||r) = 0$ iff $p = r$ and $I_q(p||r) > 0$ for $q > 0$. ■

Proposition 5 (Idempotence): Given a prior r and constraint set \mathcal{C} , let p be the posterior obtained, then $\arg \min_{u \in \mathcal{C}} I_q(u||p) = p$, i.e., taking the same information

into account twice has the same effect as taking it into account once.

Proof: This is a simple corollary of proposition 4, since $p \in \mathcal{C}$ the posterior obtained by taking p as prior and \mathcal{C} as constraint, will also be p . ■

Proposition 6 (Invariance): Given a prior r consider the constraint sets \mathcal{C}_1 and \mathcal{C}_2 , let $p = \arg \min_{u \in \mathcal{C}_1} I_q(u||r)$, then following relations hold

$$p = \arg \min_{u \in \mathcal{C}_1 \wedge \mathcal{C}_2} I_q(u||r) \quad (13)$$

$$= \arg \min_{u \in \mathcal{C}_1 \wedge \mathcal{C}_2} I_q(u||p) \quad (14)$$

$$= \arg \min_{u \in \mathcal{C}_2} I_q(u||p) . \quad (15)$$

Proof: $p \in \mathcal{C}_1$ and $p \in \mathcal{C}_2$ hence $p \in \mathcal{C}_1 \wedge \mathcal{C}_2$ so from proposition 4 both, (14) and (15) follow. We know that $p = \arg \min_{u \in \mathcal{C}_1} I_q(u||r)$ and $p \in \mathcal{C}_1 \wedge \mathcal{C}_2$ from the above two, (15) follows. ■

The result shows that if the posterior obtained from \mathcal{C}_1 is an element of \mathcal{C}_2 then applying \mathcal{C}_2 on the posterior in different ways does not result in any change.

Proposition 7 (Weak Subset Independence):

Let S_1, S_2, \dots, S_n be a partition of \mathcal{X} . Let the new information \mathcal{C} comprise about each of the conditional densities $p(x/x \in s_i)$, $i = 1 \dots n$. Thus, $\mathcal{C} = \mathcal{C}_1 \wedge \mathcal{C}_2 \wedge \dots \wedge \mathcal{C}_n$, where \mathcal{C}_i is the constraint set on the conditional densities of S_i . Then given the prior distribution r

$$p_{\mathcal{C}}^{min}(x/x \in S_i) = \arg \min_{p \in \mathcal{C}_i} I_q(p||r_i), \quad q \in (0, 1) , \quad (16)$$

and

$$\begin{aligned} I_q(p_{\mathcal{C}}^{min}||r) &= \sum_{i=1}^n u_i I_q(p_i||r_i) - \sum_{i=1}^n u_i \ln_q \frac{s_i}{u_i} \\ &+ (1-q) \sum_{i=1}^n \left(u_i \ln_q \frac{s_i}{u_i} I_q(p_i||r_i) \right) , \end{aligned} \quad (17)$$

hold where

$$p_{\mathcal{C}}^{min} = \arg \min_{p \in \mathcal{C}} I_q(p||r) ,$$

$$p_i(x) = p_{\mathcal{C}}^{min}(x/x \in S_i) ,$$

$$r_i(x) = r(x/x \in S_i) .$$

s_i are the prior probability of being in each subset, given by $s_i = \sum_{x \in S_i} r(x)$, and u_i are the posterior probability of being in each subset, given by $u_i = \sum_{x \in S_i} p_{\mathcal{C}}^{min}(x)$.

Proof: Let \mathcal{R} be the information defined by the constraint $\sum_{x \in S_i} p(x) = u_i$, then it follows from proposition 6 that

$$\arg \min_{p \in \mathcal{C}} I_q(p||r) = \arg \min_{p \in \mathcal{C} \wedge \mathcal{R}} I_q(p||r) .$$

Now we can apply proposition 2 to get (16) and (17). ■

This result is same as proposition 2 and has the same interpretation. This difference here lies in the fact that we do not have a prior information \mathcal{M} regarding the total probability in each subset.

Proposition 8 (Subset Aggregation): Let S_1, S_2, \dots, S_n be a partition of \mathcal{X} . Let Γ be a transformation which converts a given distribution p to discrete distribution over S_i , the transformation is defined by

$$p'(x_i) = \Gamma p = \int_{S_i} p(x) dx ,$$

where x_i is a discrete state corresponding to $x \in S_i$. Let \mathcal{C}' be the new information about the distribution Γp . Then for a given prior r , then

$$r(x/x \in S_i) = p_{min}(x/x \in S_i) , \quad (18)$$

$$\Gamma p_{min} = \Gamma(p_{min}) , \quad (19)$$

$$\text{and } I_q(\Gamma p_{min} || \Gamma r) = I_q(p_{min} || r) , \quad (20)$$

where $p_{min} = \arg \min_{p \in \Gamma^{-1}(\mathcal{C}')} I_q(p || r)$.

Proof: The constraint set \mathcal{C}' is defined by a set of expectations

$$\sum_{i=1}^n p'(x_i) u_m(x_i) = \langle u_m \rangle \quad m = 1 \dots M .$$

In terms of $p = \Gamma p'$ the constraint set can be represented as

$$\int_{\mathcal{X}} p(x) w_m(x) = \langle u_m \rangle \quad m = 1 \dots M ,$$

where w_m is defined as

$$w_m(x) = u_m(x_i), \quad \text{for } x \in S_i, \quad i = 1 \dots n ,$$

i.e., w_m is constant in each of the subsets S_i .

From (11) we get

$$p_{min}(x) = \frac{r(x)}{\hat{Z} \exp_q \left(- \sum_{m=1}^M \beta_m w_m(x) \right)} . \quad (21)$$

Since w_m is a constant within each subset S_i and i is a constant in itself. So equation (21) reduces to:

$$p_{min}(x) = K_i r(x) ,$$

where K_i is a constant for each subset. Now we have

$$\begin{aligned} r(x/x \in S_i) &= r(x) / \int_{y \in S_i} r(y) \\ &= p_{min}(x/x \in S_i) . \end{aligned}$$

This proves (18).

Now consider the relation

$$\begin{aligned} I_q(p_{min} || r) &= \sum_{i=1}^n u_i I_q(p_i || r_i) - \sum_{i=1}^n u_i \ln_q \frac{s_i}{u_i} \\ &\quad + (1-q) \sum_{i=1}^n \left(u_i \ln_q \frac{s_i}{u_i} I_q(p_i || r_i) \right) , \end{aligned} \quad (22)$$

which follows from (17). where

$$p_i(x) = p_{min}(x/x \in S_i) ,$$

$$r_i(x) = r(x/x \in S_i) ,$$

$$s_i = \sum_{x \in S_i} r(x) ,$$

$$\text{and } u_i = \sum_{x \in S_i} p_{min}(x) .$$

From (18) we have that $p_i(x) = r_i(x)$ and hence $I_q(p_i || r_i) = 0$. Now equation (22) reduces to

$$\begin{aligned} I_q(p_{min} || r) &= - \sum_{i=1}^n u_i \ln_q \frac{s_i}{u_i} \\ &= I_q(\Gamma p_{min} || \Gamma r) . \end{aligned}$$

This proves (19) and (20). ■

VI. SOME OBSERVATIONS ON DUALITY AND PYTHAGORAS

A. Pythagorean Property

Because of its extensive use in many problems, Pythagorean property is very important. It has been shown to exist for both second and third formalisms, involving q -expectation and normalized q -expectation respectively. In this section we have attempted to find the equivalent result for the classical expectation. The result we got is not promising but we present it here for future reference, and to introduce an alternative way to manipulate the Lagrange multipliers. Lets formally state our problem at hand:

a) Problem statement : Let r be the prior distribution and let p be the posterior got by minimizing the Tsallis divergence subject to the constraint set \mathcal{C}

$$\sum_{x \in \mathcal{X}} p(x) u_m(x) = \langle u_m \rangle \quad m = 1 \dots M .$$

Let l be another distribution satisfying the constraint

$$\sum_{x \in \mathcal{X}} l(x) u_m(x) = \langle w_m \rangle \quad m = 1 \dots M .$$

We are interested in finding the relation between $\langle u_m \rangle$ and $\langle w_m \rangle$ so as to minimize the divergence $I_q(l || p)$.

b) *Solution:* To find a solution to this problem we shall minimize the Tsallis divergence in a different manner. We start the minimization with the following Lagrangian

$\mathcal{L} =$

$$\sum_{x \in \mathcal{X}} p(x) \frac{\left[\frac{p(x)}{r(x)} \right]^{q-1} - 1}{q-1} - (1-q\lambda) (\sum_{x \in \mathcal{X}} p(x) - 1) + \sum_{m=1}^M q\beta_m (\sum_{x \in \mathcal{X}} u_m(x)p(x) - \langle u_m \rangle) ,$$

differentiating \mathcal{L} with respect to $p(x)$ and equating to 0, we get

$$\ln_q \left(\frac{r(x)}{p(x)} \right) = \lambda - \sum_{m=1}^M \beta_m u_m(x) , \quad (23)$$

$$p_{min} = p(x) = \frac{r(x)}{\exp_q(\lambda - \sum_{m=1}^M \beta_m u_m(x))} . \quad (24)$$

Multiplying equation (23) by $p(x)$ and summing it over \mathcal{X} we get

$$\begin{aligned} \sum_{x \in \mathcal{X}} p(x) \ln_q \left(\frac{r(x)}{p(x)} \right) &= \sum_{x \in \mathcal{X}} p(x) \lambda \\ &\quad - \sum_{x \in \mathcal{X}} \sum_{m=1}^M p(x) \beta_m u_m(x) , \\ -I_q^{min}(p||r) &= \lambda - \sum_{m=1}^M \beta_m \langle u_m \rangle . \end{aligned}$$

Differentiating $I_q^{min}(p||r)$ with respect to $\langle u_m \rangle$ we get

$$\frac{\partial I_q^{min}}{\partial \langle u_m \rangle} = \beta_m . \quad (25)$$

Substituting

$$\beta_m = \beta'_m (1 + (1-q)\lambda) , \quad (26)$$

equation (23) reduces to

$$\begin{aligned} \ln_q \left(\frac{r(x)}{p(x)} \right) &= \lambda \oplus_q - \sum_{m=1}^M \beta_m u_m(x) \\ p(x) &= \frac{r(x)}{\widehat{Z} \exp_q(-\sum_{m=1}^M \beta_m u_m(x))} , \end{aligned}$$

where $\widehat{Z} = \exp_q(\lambda)$. Hence equation (23) can be rewritten as

$$\ln_q \left(\frac{r(x)}{p(x)} \right) = \ln_q \widehat{Z} - \sum_{m=1}^M \beta_m u_m(x) .$$

Multiplying this equation $p(x)$ and summing it over \mathcal{X} we get

$$-I_q^{min} = \ln_q \widehat{Z} - \sum_{m=1}^M \beta_m \langle u_m \rangle .$$

Differentiating I_q^{min} with respect to β_m and equating to 0 we get

$$\frac{\partial \ln_q \widehat{Z}}{\partial \beta_m} = \langle u_m \rangle . \quad (27)$$

Equations (25) and (27) are the Legendre transform relations. Given the relations and the divergence minimization let us look at the Pythagorean property.

We want to minimize the divergence $I_q(l||p)$. For this we will proceed as follows

$$I_q(l||r) - I_q(l||p) = - \sum_{x \in \mathcal{X}} l(x) \left[\ln_q \frac{r(x)}{l(x)} - \ln_q \frac{p(x)}{l(x)} \right] ,$$

using the relation $\ln_q \left(\frac{x}{y} \right) = y^{q-1} (\ln_q x - \ln_q y)$, we get

$$\begin{aligned} I_q(l||r) - I_q(l||p) &= - \sum_{x \in \mathcal{X}} l(x) \left[\ln_q \frac{r(x)}{p(x)} \left[1 + (1-q) \ln_q \frac{p(x)}{l(x)} \right] \right] , \end{aligned}$$

using equation (23)

$$\begin{aligned} I_q(l||r) - I_q(l||p) &= - \sum_{x \in \mathcal{X}} l(x) \left[\left(\lambda - \sum_{m=1}^M \beta_m u_m(x) \right) \left(1 + (1-q) \ln_q \frac{p(x)}{l(x)} \right) \right] \\ &= \lambda - \sum_{m=1}^M \beta_m \langle w_m \rangle - (1-q)\lambda I_q(l||p) \\ &\quad - (1-q) \sum_{x \in \mathcal{X}} \left(l(x) \ln_q \frac{p(x)}{l(x)} \sum_{m=1}^M \beta_m u_m(x) \right) . \end{aligned} \quad (28)$$

The minimum of $I_q(l||p)$ is achieved for

$$\frac{\partial I_q(l||p)}{\partial \beta_m} = 0 .$$

Differentiating (28) we get

$$\begin{aligned} \frac{\partial \lambda}{\partial \beta_m} - \langle w_m \rangle - (1-q) I_q(l||p) \frac{\partial \lambda}{\partial \beta_m} \\ - (1-q) \sum_{x \in \mathcal{X}} l(x) \frac{\partial}{\partial \beta_m} \left[\ln_q \frac{p(x)}{l(x)} \sum_{m=1}^M \beta_m u_m(x) \right] = 0 . \end{aligned}$$

Using equation (27) we get

$$\begin{aligned}
\langle u_m \rangle - \langle w_m \rangle - \langle u_m \rangle I_q(l||p) \\
= (1-q) \sum_{x \in \mathcal{X}} l(x) \frac{\partial}{\partial \beta_m} \left[\ln_q \frac{p(x)}{l(x)} \sum_{m=1}^M \beta_m u_m(x) \right] \\
= (1-q) \sum_{x \in \mathcal{X}} l(x) \left[\beta_m \ln_q \frac{p(x)}{l(x)} \right. \\
\left. + \sum_{m=1}^M \beta_m u_m(x) \frac{\partial}{\partial \beta_m} \left(\ln_q \frac{p(x)}{l(x)} \right) \right] . \quad (29)
\end{aligned}$$

Evaluating it further we by using the relations $\ln_q \left(\frac{x}{y} \right) = \frac{\ln_q x - \ln_q y}{1 + (1-q) \ln_q y}$ and $p(x) = \frac{r(x)}{\lambda - \sum_{m=1}^M \beta_m u_m(x)}$, We get

$$\begin{aligned}
\langle w_m \rangle = \\
\langle u_m \rangle (1 - I_q(l||p)) + (1-q) \beta_m I_q(l||p) \\
- \sum_{x \in \mathcal{X}} \frac{l^q(x) (u_m(x) - \langle u_m \rangle) \Psi}{\left[1 + (1-q) \ln_q \frac{r(x)}{p(x)} \right]^2} , \quad (30)
\end{aligned}$$

where $\Psi = (1 + (1-q) \ln_q r(x)) \sum_{m=1}^M \beta_m u_m(x)$. Note that in this expression $r(x)$ can be replaced in terms of $p(x)$.

Though this relation does not seem promising, we have mentioned it here for the sake of completion.

B. Additive transformation - $q \leftrightarrow 2-q$

In q -deformed algebra there exists a $q \leftrightarrow 2-q$ duality. Which is the following:

$$\ln_q(1/x) = \ln_{2-q}(x) , \quad (31)$$

$$\exp_q(-x) = \frac{1}{\exp_{2-q}(x)} . \quad (32)$$

Using this duality Tsallis entropy has been well studied, i.e., various properties of S_{2-q} has been studied. Initial observations regarding S_{2-q} were made by Baldovin and Robledo [21]. Naudts [22] has further analyzed both the dualities. More study has been carried forward by Wada and Scarfone [23]. they have found relations between the Lagrange multipliers of both the dualities. In this section we introduce a similar transformation for Tsallis divergence.

Given a prior r and the constraints set \mathcal{C} defined by

$$\begin{aligned}
\sum_{x \in \mathcal{X}} p(x) &= 1 , \\
p(x) &\geq 0 , \\
\sum_{x \in \mathcal{X}} u_m(x) p(x) &= \langle u_m \rangle , \quad m = 1, \dots, M .
\end{aligned}$$

from equation (11) we have

$$p(x) = \frac{r(x)}{\widehat{Z} \exp_q \left(- \sum_{m=1}^M \beta_m u_m(x) \right)} ,$$

and using the relation (32) it becomes

$$p(x) = \frac{r(x) \exp_{2-q} \left(\sum_{m=1}^M \beta_m u_m(x) \right)}{\widehat{Z}} .$$

This form for the posterior is very good and is the basis for the $q \leftrightarrow 2-q$ transformation. Note that

$$2 - (2 - q) = q ,$$

i.e if we minimize $I_{2-q}(p||r)$ instead of $I_q(p||r)$, we have.

$$\begin{aligned}
p(x) &= \arg \min_{p \in \mathcal{C}} I_{2-q}(p||r) \\
&= \frac{r(x) \exp_q \left(\sum_{m=1}^M \beta_m u_m(x) \right)}{\widehat{Z}} .
\end{aligned}$$

VII. CONCLUSION

In this work we explored Shore and Johnson properties for Tsallis formalism of the third kind involving normalized q -expectation, it was observed that none of these properties hold for the formalism. Whereas in the study of first formalism involving classical expectation, we have been able to establish substantial number of Shore and Johnson properties. We were also able to establish a crude form of Pythagorean relation. We have also been found a $q \leftrightarrow 2-q$ additive transformation, which gives a very good form for the posterior distribution. We conclude from these observations that the first formalism is of stronger theoretical and practical significance; and these results along with the $q \leftrightarrow 2-q$ additive transformation also provides some ground work for definition of a power law family.

REFERENCES

- [1] T. M. Cover and J. A. Thomas. *Elements of Information Theory*. Wiley, New York, 1991.
- [2] Alfred Rényi. On the foundations of information theory. *Rev. Inst. Internat. Stat.*, 33:1–14, 1965. (reprinted in [24], pp. 304–317).
- [3] J. Havrda and F. Charvát. Quantification method of classification process: Concept of structural α -entropy. *Kybernetika*, 3:30–35, 1967.
- [4] Wada Tatsuaki and Saito Takeshi. When non-extensive entropy becomes extensive. *Physica A*, 301:284–290, 2001.
- [5] Hiroki Suyari. Generalization of Shannon-Khinchin axioms to nonextensive systems and the uniqueness theorem for the nonextensive entropy. *IEEE Transactions on Information Theory*, 50(8):1783–1787, 2004.
- [6] Shigeru Furuichi. On uniqueness theorem for Tsallis entropy and Tsallis relative entropy. *IEEE Transactions on Information Theory*, 51(10):3638–3645, 2005.
- [7] I. Csiszár and P. Shields. *Information Theory and Statistics: A Tutorial*, volume 1 of *Foundations and Trends in Communications and Information Theory*. Now Publications, 2004.

- [8] D.R. Clark and C.A. Thayer. A primer on the exponential family of distributions. In *Casualty Actuarial Society Spring Forum*, pages 117–148, 2004.
- [9] C.M. Bishop and SpringerLink (Online service). *Pattern recognition and machine learning*, volume 4. Springer New York, 2006.
- [10] S. Kullback. *Information Theory and Statistics*. Wiley, New York, 1959.
- [11] S. M. Ali and S. D. Silvey. A general class of coefficients of divergence of one distribution from another. *J. Roy. Statist. Soc. Ser. B*, 28:131–142, 1966.
- [12] Y. Qiao and N. Minematsu. A study on invariance of divergence and its application to speech recognition. *Signal Processing, IEEE Transactions on*, 58(7):3884–3890, 2010.
- [13] Ernesto P. Borges. A possible deformed algebra and calculus inspired in nonextensive thermostatics. *Physica A*, 340:95–101, 2004.
- [14] Constantino Tsallis, Renio S. Mendes, and A. R. Plastino. The role of constraints within generalized nonextensive statistics. *Physica A*, 261:534–554, 1998.
- [15] Lisa Borland, Angel R. Plastino, and Constantino Tsallis. Information gain within nonextensive thermostatics. *Journal of Mathematical Physics*, 39(12):6490–6501, 1998.
- [16] A. Dukupati, M. N. Murty, and S. Bhatnagar. Nonextensive triangle equality and other properties of Tsallis relative-entropy minimization. *Physica A*, 361:124–138, 2006.
- [17] A. Dukupati. Nonextensive Pythagoras’ theorem. *arXiv:cs.IT/0611030*, 2007.
- [18] A. Dukupati, M. N. Murty, and S. Bhatnagar. Properties of Kullback-Leibler cross-entropy minimization in nonextensive framework. In *Proceedings of IEEE International Symposium on Information Theory (ISIT)*, pages 2374–2378. IEEE Press, 2005.
- [19] A. Ohara. Geometry of distributions associated with tsallis statistics and properties of relative entropy minimization. *Physics Letters A*, 370(3-4):184–193, 2007.
- [20] John E. Shore. Properties of cross-entropy minimization. *IEEE Transactions on Information Theory*, IT-27(4):472–482, 1981.
- [21] F. Baldovin and A. Robledo. Nonextensive pesin identity: Exact renormalization group analytical results for the dynamics at the edge of chaos of the logistic map. *Physical Review E*, 69(4):045202, 2004.
- [22] Jan Naudts. Generalized thermostatics and mean-field theory. *Physica A*, 332:279–300, 2004.
- [23] T. Wada and A. M. Scarfone. Connections between Tsallis’ formalism employing the standard linear average energy and ones employing the normalized q -average enery. *Physics Letters A*, 335:351–362, 2005.
- [24] Pál Turán, editor. *Selected Papers of Alfréd Rényi*. Akademia Kiado, Budapest, 1976.